Dirk Pijpops* and Freek Van de Velde

# Constructional contamination: How does it work and how do we measure it?

**Abstract:** In this article, we introduce the effect of "constructional contamination". In constructional contamination, a subset of the instances of a target construction deviate in their realization, due to a superficial resemblance they share with instances of a contaminating construction. We claim that this contaminating effect bears testimony to the hypothesis that language users do not always execute a full parse while interpreting and producing sentences. Instead, they may rely on what has been called "shallow parsing", i. e., chunking the utterances into large, unanalyzed exemplars that may extend across constituent borders. We propose several measures to quantify constructional contamination in corpus data. To evaluate these measures, the Dutch partitive genitive is taken under scrutiny as a target construction of constructional contamination. In this case study, it is shown that neighboring constructions play a crucial role in determining the presence or absence of the -s suffix among instances of the partitive genitive. The different measures themselves, however, are not construction-specific, and can readily be used to track constructional contamination in other case studies as well.

**Keywords:** constructional contamination, shallow parsing, exemplar, partitive genitive, mixed-effects generalized linear models

# 1 Introduction

Constructional contamination is the effect whereby a subset of instances of a target construction is (stochastically) affected in its realization by a contaminating construction, because of a coincidental resemblance between the superficial

---

*Corresponding author: Dirk Pijpops,** Flanders Research Foundation (FWO); Research Unit QLVL, University of Leuven, Blijde Inkomststraat 21, P.O. Box 3308, 3000 BE-Leuven, Belgium, E-mail: dirk.pijpops@kuleuven.be
**Freek Van de Velde,** Research Unit QLVL, University of Leuven, Blijde Inkomststraat 21, P.O. Box 3308, 3000 BE-Leuven, Belgium, E-mail: freek.vandevelde@kuleuven.be

strings of instances of the target construction and a number of instances of the contaminating construction. This contaminating influence may not only affect the fringes of the target construction but may penetrate its very core, affecting even completely unambiguous instances. Still, the target construction and the contaminating construction typically do not merge and maintain their status as completely distinct constructions.[1] We claim that this contaminating effect is the natural result of "shallow parsing" and exemplar chunking (Bybee 2010; Dąbrowska 2012, 2014; Diessel 2015), and we do not need any additional theoretical assumptions to explain it.

As an example of constructional contamination, compare the sentences in (1) and (2), which show instances of two constructions that are structurally and etymologically unrelated.[2] While their strings, i. e., surface realizations, may look similar, their constituent structure does not, as is shown in the glosses.[3] Example (1) contains a partitive genitive construction *iets verkeerd* 'something wrong'; specifically, it is the variant without an *-s* suffix on the adjective which alternates with a variant where an *-s* suffix is realized, yielding an opposition between *iets verkeerd* and *iets verkeerd-s*. Conversely, in example (2), the quantifier *iets* 'something' forms an independent noun phrase that functions as a direct object, while *verkeerd* 'wrongly' is an adverb modifying the passive *geïnterpreteerd wordt* 'is interpreted'. In Dutch, adjectives used in adverbial function do not take the partitive *-s* suffix, so *verkeerds* would have been ungrammatical in (2).

(1)     Target construction: partitive genitive
         *in  begin     van  de  week*   ***iets***        ***verkeerd*** *gegeten*
         [in  beginning  of   the  week]PP [something  wrong]NP  eaten
                                                                    (#LEUV_4.sml)
         'I ate something wrong at the start of the week.'

---

1 To help conceive how such contaminating influence may take place between two otherwise clearly distinct constructions, consider the following comparison. The mere existence of the moon has a profound impact on life on earth, with many human fishing communities and the life of many animals and plants being centered around the flow of the sea tides. Despite this influence, earth and moon remain two clearly distinct planetary bodies.
2 These, as well as all other corpus examples in this paper, were taken from the *ConDiv* corpus of written Dutch (Grondelaers et al. 2000). Next to each example, the corpus file can be found from which the example was taken.
3 The bracketed notation is only added for expository reasons, and should reflect an uncontroversial constituent structure. We do not mean to subscribe to any specific theoretical framework by using these brackets.

(2)    Contaminating construction: construction with adverb
       *dat* **iets**           **verkeerd**      *geïnterpreteerd wordt?*
       that [something]$_{NP}$ [wrongly]$_{AdvP}$ interpreted    gets

                                                                (#VLAAN_1.sml)

       '...that something gets wrongly interpreted?'

While *iets* and *verkeerd* happen to occur next to each other in (2), there is no obvious reason why they should occur together in other instances of this construction. In fact, other instances of the partitive genitive and the construction with an adverb may look totally different from each other, such that the superficial resemblance between (1) and (2) is purely coincidental. However, the noun phrase *iets* and the adverb *verkeerd* do happen to occur adjacent to each other quite frequently. It is not uncommon that speakers express the state of affairs that *something* is being *wrongly* interpreted, done, or understood.

In this paper, we show that the frequent co-occurrence of the contiguous expression of the quantifier *iets* and the adverb *verkeerd* in the construction in (2) generates a measureable preference for the variant without *-s* in partitive genitives such as (1). *Iets* and *verkeerd* are just examples, of course, and we shall identify a number of other such contaminating co-occurrences. In fact, this contaminating effect is found to be the main determinant of the occurrence of this *-s* suffix, in that it vastly outperforms often cited regional and stylistic factors (van der Horst 2008: 1624–1625; Broekhuis 2013: 426).[4]

In Section 2, the effect of constructional contamination is explained in detail. Next, Section 3 presents the case study that is employed to exemplify and study constructional contamination in this article, i.e., the Dutch partitive genitive construction. Section 4 then describes the extraction of the relevant data, as well as the results of a distinctive collexeme analysis. Section 5 forms the bulk of the article. Here, we propose four quantitative measures of constructional contamination, evaluate them against the case study at hand, and explain

---

**4** In this example, and in the rest of the paper, the partitive genitive construction is put under the microscope as the target construction of constructional contamination. However, there is no reason to assume that constructional contamination is strictly unidirectional. In fact, we expect it to be typically bidirectional, with the partitive genitive in (1) also contaminating the construction with adverb in (2). There are some indications that this is indeed the case. The internet contains a fair number of instances of the construction in (2) in which a reading as a partitive genitive is infelicitous, such as *iets verkeerds geïnterpreteerd* 'interpreted something wrongly' or *iets verkeerds gelopen* 'something gone wrong', yet which still receive a – normally ungrammatical – *-s* ending. We see no other way to explain this ungrammatical *-s* ending than through constructional contamination.

what the results of this evaluation tell us about the exact nature of constructional contamination. Finally, Section 7 summarizes the conclusions.

# 2 Constructional contamination

## 2.1 Links in the construction

Many linguists view language as a structured array of conventional form–function pairings (Langacker 2008: 222), commonly called "signs" or "constructions". The "array" or "repository", in which all constructions – whether atomic or complex, schematic or concrete – are stored, is often referred to as the "constructicon" (a term supposedly coined by Jurafsky 1992). Constructional linguists (Lakoff 1987; Goldberg 1995; Croft 2001) have repeatedly pointed out that it is organized as a network: the constructicon is not like a disorderly set of discrete form–meaning pairs, much like a drawer in which pairs of socks are lying around. Rather, constructions entertain vertical relationships in which more schematic constructions subsume lower, concrete constructions that "inherit" features from the dominating nodes, as well as horizontal relationships in which constructions stand in differential opposition to each other (Van de Velde 2014).

Despite their interconnectedness, the form–meaning pairing of constructions should, at first sight, be as fixed and predictable as possible, lest the unique semiotic link between a form and its function be jeopardized, and lest the Saussurian-style horizontal opposition relations collapse. This is the reason why structuralism had a hard time explaining homonymy, synonymy, and language change, and adhered to isomorphism (see the discussion in Croft 2001: 111–119). If language is "un système où tout se tient", in which symbolic units exist by virtue of the differential opposition links they entertain with each other, shifts in the system may bring about a collapse of the system. If meaning A corresponds to forms {X, Y, Z}, and form X corresponds to meanings {A, B, C}, thus displaying a many-to-many mapping, then language users face difficulties in decoding and encoding language. The naïve structuralist conception of the constructicon with interconnected constructions that are themselves nevertheless "claires et distinctes" entails then that constructions should not be contaminated by neighboring constructions: constructions ought to be delineated from one another as discretely as possible. Fluid boundaries are semiotically problematic. Indeed, in language acquisition, it has been shown that children have difficulties acquiring the second member of a synonymy set if they have already acquired the first, suggesting that they operate with a one-form-one-

meaning heuristic (Markman and Wachtel 1988; Markman et al. 2003; Abbot-Smith and Behrens 2006).

This ideal, pure, non-promiscuous nature of the linguistic sign is, however, at variance with reality. Constructions are known to infect each other on the formal as well as on the semantic level, for instance when formal similarities between two constructions cause semantic convergence, or when semantic similarities cause distributional convergence of the forms (see De Smet 2010, on what he calls "grammatical interference"). Indeed, diachronically, a construction often derives from multiple lineages that come to merge (see Van de Velde et al. 2013 on "multiple source constructions"), and synchronically, a construction often displays contamination effects at its fringes.[5]

Examples of diachronic merger of constructions can be found in Van de Velde et al. (2013), and superficial morphological similarities that cause diachronic convergence are discussed in Van de Velde and van der Horst (2013), who use the term "homoplasy" for this phenomenon. In other cases, the integration of two clauses is less seamless, and stitches are still apparent (De Smet and Van de Velde 2013, working on the diachronic correlate of what Lakoff 1988 [1974] has called "syntactic amalgams").

As an example of synchronic contamination effects, let us consider the use of the preterite subjunctive in backshifted clauses like (3) (Huddleston 2002: 87). Varieties of English, such as American English, which prefer the subjunctive in mandative clauses, normally use the present subjunctive *be* here, see (4), reserving the past subjunctive for irrealis conditionals, as in (5), as opposed to the indicative in (6) (Bergs and Heine 2010: 110–111). In (3), there is a contaminating influence of the present vs. past opposition in the indicative on the subjunctive. This contamination is founded on a formal similarity link: in the verb *be* the past subjunctive *were* has formal ties to the past indicative *were*.

(3)  *The head of department insisted that he were promoted.*
(4)  *The head of department insisted that he be promoted.*
(5)  *If he were here, he would beg to differ.*
(6)  *If he is here, we should have the meeting right now.*

---

**5** The diachronic merger of different constructions may of course originate as synchronic contaminations at the periphery of constructions (see for instance Fonteyn and van de Pol 2016 for a detailed case study in English).

The upshot of this paper is that constructions may affect other constructions through such "superficial links", i. e., through strings that are similar at the surface level, irrespective of their structural differences, such as the link between (1) and (2). At the same time, however, there is an important role for semantics. The impact of superficially similar constructions will be shown to be stronger to the extent that semantic ambiguity is involved, that is, if a particular string can be interpreted as an instantiation of each of the two superficially resembling constructions.[6] This is reminiscent of what in diachronic grammar is known as "bridging contexts" (Heine 2002), where reanalysis is facilitated by or even entails an intermediate ambiguous stage. A case in point is the development of *be going to* from a lexical verb of "motion + intention", as in (7), to an auxiliary conveying "imminent future" in (9). The bridge between the old and the new reading are cases like (8), where the old motion reading is still possible.

(7)   *He is going to the market to buy oranges.*
(8)   *He is going to tell her the truth.*
(9)   *He is going to wake up any minute now.*

In the case of constructional contamination, however, no reanalysis is involved. Rather, superficially similar but grammatically independent structures exert an influence on the formal realization of the target construction, and this effect is not only seen in bridging cases, but even in completely unambiguous instances, as will be shown below.

---

**6** This may sound complicated, and a non-linguistic comparison may be of help here. Men's first flutes were made from animal bones with holes drilled in them. Now suppose, for the sake of the argument – the actual history is not at issue here (interested readers may be referred to http://cogweb.ucla.edu/ep/FluteDebate.html) – that our ancestors hit upon that idea because animal bones with holes were already fortuitously available in the form of prey animals bones pierced by carnivorous animals' teeth. What we would argue under the scenario of constructional contamination is that the man-crafted flute design would be partially influenced by earlier, naturally occurring punctured bones. The interspacing of the holes could conform to the distance of the carnivorous teeth, for instance, even if that interspacing was not optimal for the range of tones that aesthetically appeals to human ears. The "sensitivity to semantics" in constructional contamination could be understood as the influence on the flute design only being apparent when men actually use bone material to manufacture flutes. The design of flutes made from, say, tree branches, would not be influenced by naturally occurring punctured bones, and the interspacing of the holes would be geared towards optimal auditory aesthetics.
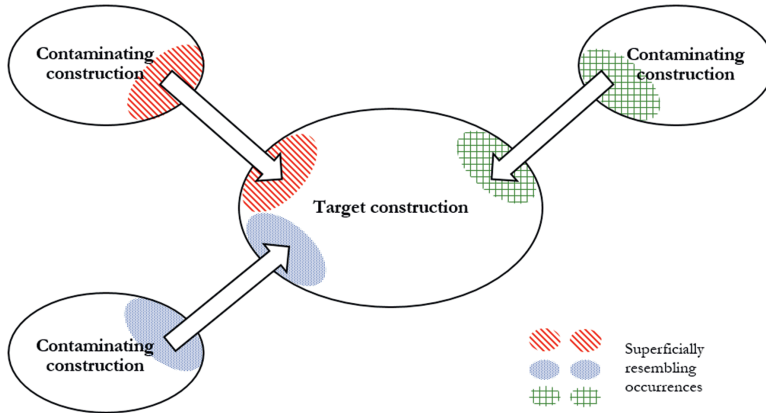
**Figure 1:** Constructional contamination through superficially resembling occurrences of three contaminating constructions and a target construction.

The mechanism of constructional contamination is represented schematically in Figure 1. It visualizes the contaminating influence from various neighboring constructions on a target construction, bridging the constructions' borders. One pillar of this bridge is formed by occurrences of the neighboring constructions which superficially resemble the target construction, as in (2); the other pillar by instances of the target construction that superficially resemble these other constructions, as in (1).

If constructions indeed contaminate one another, this is a clear infringement on the principle of isomorphism and messes up the neat "bijection" (to borrow a term from mathematics) between form and function. Going back to the idea of the "constructicon" as a network, we can then say that the interconnectedness of constructions goes beyond the vertical inheritance relations and the horizontal opposition relations. Constructions are also intertwined horizontally by resemblance relations (see also Norde 2014), whereby semantic affinity may beget formal affinity and vice versa. The constructicon network is thus rather like a tangled ball of wool, where individual threads cannot be isolated easily, as they are inextricably entangled in other threads, and several may even get interwoven into a single thread.

This may all sound fine, but how should this contaminating interconnectedness between constructions be conceived of concretely? This question is answered in two ways. On the theoretical level, Subsection 2.2 explains what constructional contamination tells us about how language users process language. It is claimed that the effect of constructional contamination emerges as

a side-effect of exemplar chunking. On the methodological level, in Sections 4 and 5, we will use distinctive collexeme analysis and especially mixed-effects logistic regression (Baayen 2008; Gries 2013a; Speelman 2014) to measure the impact of adjacent constructions.

## 2.2 Chunking

Taken in isolation, (1) and (2) are two completely normal Dutch sentences. Languages users should have no difficulties in parsing them, and there is no a priori reason why (1) and (2) should interfere with each other. Their resemblance is only superficial and coincidental. Still, one may doubt whether language users actually carry out a full syntactic parse of the strings they hear or pronounce. Fully analyzing syntactic structures is a computationally demanding task, especially under time pressure of online language processing. As an alternative, it has been observed in earlier studies that a "pseudo-parse" is often sufficient for communicative success (see Dąbrowska 2012, referring to Townsend and Bever 2001, Ferreira and Patson 2007, and other publications). In such a pseudo-parse, language users do not execute a full parse of a sentence, but rather make use of short-cuts like exemplar chunking for more efficient processing.

Instead of fully analyzing (2) to its underlying grammatical structure, the language user may simply store *iets verkeerd* as a chunk. Frequent co-occurrence of these two elements can lead to a bond that overrides constituent borders (Bybee 2010). Later, when this same language user needs to express something superficially similar, like the sentence in (1), he/she can easily access this ready-made *iets verkeerd*, without *-s* suffix, instead of having to compose it from scratch. This strategy increases processing efficiency, at the cost of memory usage: while *iets verkeerd* no longer needs to be fully rebuilt or parsed, it does need to be stored in memory.[7]

For the current article, the crucial claim here is that the language user does not care whether the underlying grammatical structure of *iets verkeerd* in (1) and (2) is different or the same. All is well as long as both syntactic strings "look" similar enough in semantics and form, i. e., when there is superficial semantic and formal resemblance. Therefore, the more frequently language users hear instances such as (2), the stronger *iets verkeerd* gets entrenched, and the more often *-s* drop can be expected in superficially resembling partitive genitives. This

---

**7** Although human memory is of course finite, it is fairly large (see Bartol et al. 2015 and references cited therein for estimates).

contagious effect of the non-partitives on -s drop in genuine partitives thus goes to show that language users are sensitive to frequency effects of exemplars of a particular multiword string (see also Arnon and Snider 2010).[8]

Note that this means that the contaminating links described in Subsection 2.1 are only real at two levels. The first is the level of language description; this is simply to say that the effects of constructional contamination can be assessed in corpus data (see Section 5). The second is the procedural level of language processing. By this, we mean that the cause of the contamination effects lies in the short-cuts that language users may take while parsing or producing utterances, as explained above. Perhaps future work may claim that contamination links are also cognitively real at the declarative level. That is to say that these links exist in people's minds even when they are not processing language. However, at the moment, this is an assumption that is not needed to explain the effects we are observing, and we will therefore make no claims in this direction.

# 3 Dutch partitive genitives

The partitive genitive of Dutch is a relic construction. In Present-day Dutch, partitives are expressed as close-apposition binominals without genitive inflection, which they still had in Middle or Early Modern Dutch; see (10) vs. (11) and (12) vs. (13).[9] There exists, however, one inconspicuous corner in contemporary Dutch grammar where a partitive genitive inflection survives, namely in contexts where an indefinite pronoun or quantifier expression is followed by an adjective (Haeseryn et al. 1997: 863; Booij 2010: 223–228; Broekhuis 2013: 420–426); see example (14).[10] Without attributing too much importance to the

---

**8** Constructional contamination is furthermore related to what Szmrecsanyi has called "β-persistence", i. e., a priming effect of patterns that are not real instantiations of the target, but share formal characteristics with it, for instance, a non-comparative use of *more* (e. g., *This movie has more violence in it*) acting as a prime for an analytic comparative (*more interesting*) (Szmrecsanyi 2005: 140). The difference is that in cases of β-persistence, the frequency in incremental discourse is affected when the prime or contaminator is in the vicinity of the target, whereas in cases of constructional contamination, we have an effect across the board on the target. Put in other terms: β-persistence is a syntagmatic effect, whereas constructional contamination is more paradigmatic in nature.

**9** There are a number of other competing constructions as well in Present-day Dutch, which we will gloss over here. See Van de Velde (2009: Ch.3) and Hoeksema (2014) and references cited there, for details.

**10** We will not be concerned with the terminological distinction between indefinite pronoun, quantifier, indefinite numeral, etc., which is made in reference grammars of Dutch (Haeseryn et al. 1997: 432; van Bart et al. 1998: 17–28; Balk-Smit Duyzentkunst 2000: 78–103).

formalism itself, we can represent the constructional template as in (15) (following Booij 2010: 227): the semantic pole is on the right-hand side of the double arrow and the formal pole on the left-hand side, and $NP_i$ represents a quantifier expression.

(10)　*een pont　speck-s*
　　　a　pound　bacon-GEN
　　　'a pound of bacon'
　　　(Middle Dutch, van der Horst 2008: 575)

(11)　*een pond　spek*
　　　a　pound　bacon
　　　'a pound of bacon'

(12)　*een corste　broot-s*
　　　a　crust　bread-GEN
　　　'a crust of bread'
　　　(Early Modern Dutch 16th century, van der Horst 2008: 1033)

(13)　*een korst　brood*
　　　a　crust　bread
　　　'a crust of bread'

(14)　*iets　　　bijzonder-s*
　　　something　special-GEN
　　　'something special'

(15)　$[NP_i [... [X\text{-}s]_A]_{APj}]_{NPk} \leftrightarrow [Quantity_i \text{ with } Property_j]_k$

As mentioned in Section 1, the partitive construction in (14) displays variation. It may have two realizations: one is the [+s] variant, in which an -*s* ending is added to the adjective, as in (17); the other is the [+∅] variant, in which the bare adjective is used, as in (18).[11] There does not seem to be any noticeable

---

**11** Note that this -*s* ending is a property of the construction, not of the adjective, as opposed to, for instance, the adjectival ending in German *etwas Tolles* 'something fun'. Here, the ending is part of normal German adjectival inflection, and changes as the case of the phrase changes. In Dutch however, this -*s* ending is independent of case and may always and only be applied if the adjective is part of a partitive genitive construction.

difference in semantics between both variants. The variation can be represented by putting the -s in [X-s] between round brackets, as in (16).

(16)   $[NP_i [... [X(-s)]_A]_{APj}]_{NPk} \leftrightarrow [Quantity_i$ with $Property_j]_k$

(17)   [+s] variant
    *Is  er    nog  **iets**    **leuk-s**  te  beleven?*    (#holl_6.sml)
    is  there  still  something  fun-GEN  to  experience?
    'Is there still something fun to do?'[12]

(18)   [+∅] variant
    *of      er    hier  nog  **iets**    **leuk** te  beleven    valt*
    whether  there  here  Still  something  fun   to  experience  falls
                                                (#holl_6.sml)
    '... whether there is still something fun to do here?'

While the two variants do not show any observable semantic differerence, Pijpops and Van de Velde (2014) found that the addition of the -s is probabilistically determined by a number of language-internal and language-external factors, which are listed in Table 1. Overall, the [+s] variant is most frequently used, yet the [+∅] variant is also fairly common. The variation in the realization of the partitive genitive -s offers an interesting opportunity to investigate

**Table 1:** Factors determining -s omission with the partitive genitive in Dutch.

| Variable | Influence |
| --- | --- |
| *Type-Adjective* | Increased [+∅] with the adjectives *verkeerd* 'wrong', *goed* 'good', *beter* 'better', and *fout* 'incorrect', and with the color adjectives |
| *Variety & Quantifier* | Increased [+∅] in Belgium, but only with the quantifiers *iets* 'something' and *niets* 'nothing'. In the Netherlands, there is no differentiation between the quantifiers. |
| *Register* | Increased [+∅] in the more informal registers |
| *Frequency* | Increased [+∅] in low frequent phrases |

**12** In the glosses, the -s ending has been marked -GEN, because it historically descends from a genitive marker; hence the name of the partitive genitive. The actual genitive case has been extinct in Dutch for some time though, and synchronically, the -s can perhaps better be viewed as an isolated suffix (Pijpops and Van de Velde 2016).

the mechanism of constructional contamination, as concrete strings that instantiate the pattern on the left-hand side of (16) also occur in syntactically different constructions such as (2).

# 4 Data

## 4.1 Extraction

The analysis presented here builds on the data gathered in Pijpops and Van de Velde (2014). We will adopt the dataset from this earlier study, subject it to new techniques, manually and automatically add new variables to it, and finally analyze this new information (see Sections 4 and 5). As a preliminary to our analysis, the current subsection shortly summarizes the extraction and design of the dataset, as described in Pijpops and Van de Velde (2014: 8–14).[13] The aim of that study was to assess the factors that exert an influence on the partitive -s realization, and its results are summarized in Table 1. All analyses mentioned in the present subsection were executed in the study of Pijpops and Van de Velde (2014). All analyses mentioned in the following (sub)sections were executed in the current study.

As our data source, we employed the synchronic component of the *ConDiv* corpus of written Dutch (Grondelaers et al. 2000). The *ConDiv* corpus comprises material from the Netherlands and Belgium, and is stratified according to register, containing chat logs, e-mails, mass newspapers, and quality newspapers from around the turn of the current century, totaling about 45 million words.[14] As such, it provides a representative cross-cut of Dutch written language.

The quantifier + adjective strings extracted from the *ConDiv* corpus met a number of criteria. As for the quantifiers, they had to be listed as indefinite pronouns or numerals in Haeseryn et al. (1997: 356, 432); this excluded more complex quantifier expressions like *het weinige opwindends* lit. 'the little exciting' (Booij 2010: 228), which are exceedingly scarce anyway. In addition, they had to occur in the *Corpus of Spoken Dutch* (*Corpus Gesproken Nederlands*, CGN) with at least 14 occurrences in which they are part of a partitive genitive

---

**13** We made use of AntConc (Anthony 2011) for the extraction of the corpus data.
**14** The corpus also contains the Bulletins of Acts, Orders and Decrees; they were not used as they did not provide enough occurrences to reliably fit this register in the regression models.

(Oostdijk et al. 2002).[15] Still, *iemand* 'someone' and *niemand* 'no one' were not selected, because combinations with these quantifiers do not originate from real partitive genitives, but are analogically calqued from combinations with *iets* (see WNT, s.v. *ander* and Van de Velde 2009: 107). As for the adjectives, first, they had to yield at least 7 occurrences in the *Corpus of Spoken Dutch*, in which they formed a partitive genitive together with one of the selected quantifiers. Second, the bare form of the adjective could not end on *-s* or *-isch*, as there is no graphemic or phonological differentiation between its [+∅] and [+s] variants. Third, the adjective could in no form be homographic with the plural form of a noun, as in *ouders* ('older.GEN' or 'parent.PL') or *extra's* ('extra.GEN' or 'bonus.PL'). To this list, we added the main color adjectives, as well as the adjective *beter* 'better', because they were of special interest (Van de Velde 2001: 150–151; Pijpops and Van de Velde 2014: 10). These criteria yielded the following quantifiers and adjectives.

> Quantifiers: *iets* 'something', *niets* 'nothing', *wat* 'something', *veel* 'a lot', *weinig* 'few', *zoveel* 'so much'

> Adjectives: *belangrijk* 'important', *beter* 'better', *bijzonder* 'particular', *blauw* 'blue' *boeiend* 'fascinating', *concreet* 'concrete', *deftig* 'distinguished', *dergelijk* 'similar', *erg* 'awful', *geel* 'yellow', *gek* 'crazy', *goed* 'good', *grappig* 'funny', *groen* 'green', *interessant* 'interesting', *lekker* 'tasty', *leuk* 'fun', *lief* 'sweet', *mooi* 'beautiful', *nieuw* 'new', *nuttig* 'useful', *oranje* 'orange', *origineel* 'original', *positief* 'positive', *purper* 'purple', *raar* 'weird', *rood* 'red', *slecht* 'bad', *spannend* 'exciting', *speciaal* 'special', *verkeerd* 'wrong', *wit* 'white', *zinnig* 'sensible', *zwart* 'black'

All instances in which one of these quantifiers preceded one of these adjectives were extracted from the corpus. This dataset was then manually checked to exclude all false positive instances, as well as to fix those instances in which the automatic annotation failed. This left us with in total 3,018 occurrences, of which 2,388 exhibited the [+s] variant and 630 the [+∅] variant. During this manual checking, it became apparent that a disproportionally large number of the false positives contained the adjectives *verkeerd* 'wrong', *goed* 'good', *beter* 'better', and *fout* 'incorrect', as well as the color adjectives. We suspected that this could have had some effect on the instances retained with these adjectives, which could distort further analyses. To guard against this, the variable *Type-Adjective* was added. This variable distinguished between the

---

**15** We made use of the *Corpus of Spoken Dutch* for the selection of quantifiers and adjectives, because we needed an annotated corpus in which we could automatically distinguish between partitive genitives and non-partitive genitives.

color adjectives, the "assessment adjectives" *verkeerd* 'wrong', *goed* 'good', *beter* 'better', and *fout* 'incorrect', and all other adjectives.[16] Possible interactions between this variable and the other variables used in Pijpops and Van de Velde (2014) were then fed into a stepwise variable selection procedure, yet none were selected. For the detailed analysis of the regression output, we refer the reader to Pijpops and Van de Velde (2014). The variables that were shown to have an influence on the -*s* realization in partitive genitives are mentioned in Table 1.

## 4.2 Exploration

In order to allow operationalization in a regression model, the variable *Type-Adjective* generalized over several adjectives, which were grouped into the categories *color adjectives*, *assessment adjectives*, and *other adjectives*. In doing so, it is easy to lose track of the behavior of the individual adjectives.[17] As the frequency of the adjectives shows a roughly Zipfian distribution (see Figure 2, Zipf 1932), the overall behavior of a category may be strongly affected by the morphological preference of a single highly frequent adjective, while the remaining adjectives do not share this preference. This would mean that what we suspect to be the effect of constructional contamination is no more than the idiosyncratic behavior of a single adjective. The question we then need to answer is: do the color adjectives and assessment adjectives *beter, goed, fout*, and *verkeerd* fully merit their status as separately operationalized categories, or can their collective behavior be explained by the idiosyncratic preference of a single highly frequent member?

Below, we apply the exploratory technique of Distinctive Collexeme Analysis to investigate whether this is indeed the case (Gries and Stefanowitsch 2004). A Distinctive Collexeme Analysis quantifies the attraction between a lexeme and a construction or a constructional variant, in contrast to another construction

---

**16** The *Type-Adjective* variable is presented in more detail in Subsection 5.2. The term "assessment adjectives" is strictly a convenience label. Semantically, these adjectives are related, as they all express qualification (or assessment). Still, the terms "qualifier", "evaluation" or "appraisal" adjectives were not used, as these have a technical meaning in Systemic Functional Linguistics (Halliday and Matthiessen 2004) and related frameworks (Martin and White 2007 on appraisal and Hunston and Thompson 2001 on evaluation).

**17** In Pijpops and Van de Velde (2014), this was accounted for by adding the phrase as a random factor to the regression model, which did not lead to a decrease in the variable importance of the *Type-Adjective* predictor.
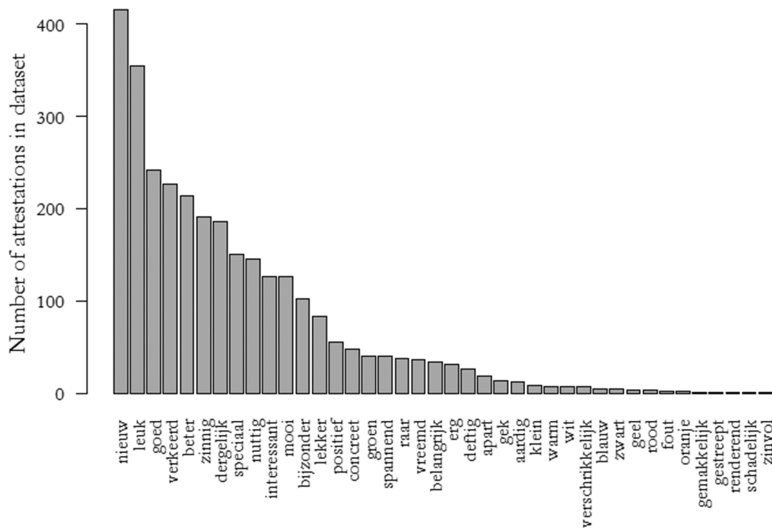
**Figure 2:** Zipfian frequency distribution of the adjectives in the dataset.

or variant.[18] In this case, the lexemes correspond to the adjectives and the constructional variants to the [+∅] and [+s] variants. The analysis indicates which adjectives are heavily attracted to the adjectival slot of the [+∅] variant of the partitive genitive construction, and vice versa, as opposed to the [+s] variant. The degree of this mutual attraction is expressed in terms of *collocational strength*.[19] Table 2 shows the results of this analysis as well as the raw frequencies per adjective. The adjectives are listed in order in decreasing collocational strength, i. e., in decreasing preference for respectively the [+∅] and [+s] variants. The assessment adjectives are shaded in dark gray, the color adjectives in light gray.

As can be seen in Table 2, the assessment and color adjectives rank the highest in terms of preference for the [+∅] variant. The only reason why *apart* 'separate' ranked higher than *fout* 'incorrect' and *oranje* 'orange' in collocational strength with the [+∅] variant is the low overall frequency of each of these

---

**18** See Gries and Stefanowitsch (2004) and Gries (2013b) for a presentation and further discussion of Distinctive Collexeme Analysis, as well as the other members of the increasingly popular family of collostructional techniques. To carry out the analysis, we made use of Gries' R-script for collocational analyses (Gries 2014).

**19** The collostructional strength measure is the negative base-10 logarithm of the *p*-value of a one-tailed Fisher-Yates Exact test on the contingency table of the adjective at issue versus all other adjectives by the occurrence in either of the two constructions.

**Table 2:** Results of the Distinctive Collexeme Analysis on the adjective. The adjectives *verkeerd, goed, beter*, and *fout* and the color adjectives all share an outspoken preference for the $[+\varnothing]$ variant.

| Preference for the $[+\varnothing]$ variant Total number of occurrences: 2388 | | | | Preference for the $[+s]$ variant Total number of occurrences: 630 | | | |
|---|---|---|---|---|---|---|---|
| **Adjective** | **$[+\varnothing]$ occ.** | **$[+s]$ occ.** | **Collostr. strength** | **Adjective** | **$[+\varnothing]$ occ.** | **$[+s]$ occ.** | **Collostr. strength** |
| *verkeerd* 'wrong' | 150 | 76 | 53.48 | *dergelijk* 'similar' | 3 | 183 | 15.18 |
| *groen* 'green' | 41 | 0 | 28.35 | *leuk* 'fun' | 23 | 331 | 14.53 |
| *goed* 'good' | 75 | 167 | 4.13 | *nieuw* 'new' | 38 | 377 | 11.15 |
| *wit* 'white' | 7 | 1 | 3.96 | *bijzonder* 'extraordinary' | 2 | 101 | 8.05 |
| *geel* 'yellow' | 4 | 0 | 2.72 | *mooi* 'beautiful' | 11 | 116 | 3.86 |
| *beter* 'better' | 62 | 152 | 2.65 | *zinnig* 'sensible' | 28 | 163 | 1.81 |
| *blauw* 'blue' | 4 | 1 | 2.10 | *lekker* 'tasty' | 10 | 73 | 1.59 |
| *zwart* 'black' | 4 | 1 | 2.10 | *gek* 'crazy' | 0 | 14 | 1.43 |
| *apart* 'separate' | 8 | 11 | 1.53 | *nuttig* 'useful' | 22 | 124 | 1.35 |
| *fout* 'incorrect' | 2 | 0 | 1.36 | *vreemd* 'weird' | 4 | 33 | 1.05 |
| *oranje* 'orange' | 2 | 0 | 1.36 | *positief* 'positive' | 8 | 47 | 0.80 |
| *deftig* 'decent' | 9 | 17 | 1.13 | *concreet* 'concrete' | 8 | 40 | 0.52 |
| *raar* 'weird' | 11 | 27 | 0.82 | *spannend* 'exciting' | 7 | 33 | 0.42 |
| *rood* 'red' | 2 | 2 | 0.71 | *klein* 'small' | 1 | 8 | 0.39 |
| *gemakkelijk* 'easy' | 1 | 0 | 0.68 | *erg* 'awful' | 6 | 25 | 0.28 |
| *warm* 'warm' | 3 | 5 | 0.65 | *aardig* 'nice' | 2 | 10 | 0.28 |
| *speciaal* 'special' | 35 | 115 | 0.60 | *verschrikkelijk* 'horrible' | 1 | 6 | 0.26 |
| *interessant* 'interesting' | 29 | 98 | 0.49 | *belangrijk* 'important' | 7 | 27 | 0.23 |
| | | | | *gestreept* 'striped' | 0 | 1 | 0.10 |

adjectives.[20] The remaining assessment adjectives, *verkeerd* 'wrong', *goed* 'good', and *beter* 'better' are all of similar frequency and all show an outspoken preference for the $[+\varnothing]$ construction. So, while *goed* and *beter* occur more frequently in the $[+s]$ variant than in the $[+\varnothing]$ variant in absolute numbers, the $[+\varnothing]$ variant is much more likely to occur here than with other adjectives. The category of assessment adjectives is thus not dominated by a single adjective that

---

**20** *Fout* 'incorrect' has so few hits in our dataset because the adjective was not included in the original set of adjectives selected for extraction (see Subsection 4.1). It only found its way in because of partitive genitive occurrences where the adjective is pre-modified by an adverb (Pijpops and Van de Velde 2014: 10–11).

explains the whole effect; rather, all these four adjectives exhibit deviant behavior in strongly preferring the $[+\varnothing]$ variant as compared to all other adjectives.

The color adjectives do have one adjective spiking in frequency, namely *groen* 'green'. This is due to a specialized use of *groen*, which is notably frequent in partitive genitive contexts. In (19), *veel groen*, literally 'a lot of green things', is used to refer to verdure.

(19)  *een   voormalige  boerderij,  omzoomd  door  kortgeschoren    hagen,*
      a     former       farm        rimmed    by     closely-clipped  hedges
      ***veel   groen*** *en   een tuin*                                 (verstr3.txt)
      much   green  and  a     garden
      'A former farmhouse, rimmed by closely-clipped hedges, a lot of verdure,
      and a garden...'

However, the other color adjectives do not behave differently from *groen* 'green'. They all exhibit a strong preference for the $[+\varnothing]$ variant. The high frequency of *groen* itself does not fully account for the effect at hand. The distinctive collexeme analysis thus confirms our intuition that the behavior of the *Type-Adjective* variable is not due to idiosyncratic behavior of one or a few hyper-frequent adjectives, but rather that all assessment adjectives and color adjectives share an outspoken preference for -*s* omission.

# 5 Measuring constructional contamination

Now that the data have been presented and it has been shown that a number of adjectives do exhibit strongly deviant preferences, we turn to the question of how constructional contamination should be measured in practice.

First, those instances of the target construction need to be identified that show some superficial resemblance to frequently occurring instances of a possible contaminating construction. In Subsection 5.1, four possible ways of identifying these instances are presented.

Second, it should be investigated whether strictly unambiguous instances of the target construction are affected. This point is important, as it would be somewhat trivial to claim that a superficially resembling neighboring construction exerts its influence on the formal realization of another construction if that other realization cannot unambiguously be classified as an instance of the target construction. For instance, (20) presents an instance that is ambiguous between a partitive genitive and a construction with an adverb, as is shown in the

translations. Of course, we expect these instances to exhibit a larger proportion of -s omission, as some of them might not actually instantiate partitive genitives, but rather constructions with adverbs, which can only appear grammatically without -s ending. Instead, what we want to show is that constructional contamination penetrates the very core of the partitive genitive and even affects strictly unambiguous partitive genitives. As such, the four measures of constructional contamination will be evaluated against a dataset only containing strictly unambiguous instances of the partitive genitive. This evaluation is described in Subsection 5.2.

Finally, a comparison of the four measures will shed light on the question just how much superficial resemblance between the contaminating constructions is needed. However, we should not just stop at comparing their performance, but should also closely examine the measures themselves to find out why some perform better than others. This is done in Subsection 5.3.

(20) *Ze    zijn   dus   vaak   erg    onzeker  en    bang   dat   ze    **iets***
     they  are   thus  often  very   insecure  and   afraid  that  they  something
     ***verkeerd** zullen  doen.*                                (verstr4.txt)
     wrong       will    do
     'As such, they are often very insecure and afraid to do something wrong.'
     **[partitive genitive]**
     'As such, they are often very insecure and afraid to do something wrongly.'
     **[construction with adverb]**

## 5.1 Measures of constructional contamination

In this subsection, four measures of constructional contamination are presented: (i) *Type-Adjective*, (ii) *Partial String Resemblance*, (iii) *String Resemblance*, and (iv) *Semantic String Resemblance*. The variable *Type-Adjective* is adopted from Pijpops and Van de Velde (2014). It is based solely on intuitive observations gained from manually checking the dataset during the study of Pijpops and Van de Velde (2014), and it is verified in the present study by a post-hoc check with a collexeme analysis (see Subsection 4.2). However, it should be possible to formulate a more direct operationalization preferably to be calculated (semi-) automatically. To achieve this, the last three measures are introduced, ordered from coarse-grained to fine-grained, with increasing superficial resemblance. Below, the four measures are presented in more detail.

*Type-Adjective* is a simple categorical variable that distinguishes between the so-called assessment adjectives *verkeerd* 'wrong', *goed* 'good', *beter* 'better',

and *fout* 'wrong', the color adjectives, and all other adjectives. So far, we have only used the adjective *verkeerd* as an example of constructional contamination. We now broaden our scope to the other assessment adjectives and the color adjectives. In what contaminating constructions do they appear?

Like *verkeerd* 'wrong', the other assessment adjectives also habitually seem to occur in constructions with adverbs that superficially resemble partitive genitives, as in (21) and (22). In addition, the same adjectives also frequently appear in another construction superficially resembling the partitive genitive, namely the predicative construction. In (23), a genuine partitive genitive with the adjective *beter* 'better' is shown. In (24), this interpretation as a partitive genitive is implausible; instead, we are dealing with a predicative construction.

(21)  *Misschien moeten we voortaan   de spelregels   toch* **iets**
      Perhaps  should  we  henceforth  the  game-rules  still  something
      **beter**  *uitleggen.*                                    (hbvl1.txt)
      better  explain
      'Perhaps from now on, we should explain the rules of the game a bit better.'

(22)  *en   ruim   van tevoren* **iets**       **goed** *plannen zodat iedereen*
      and  amply  of  before  something  well  plan   so    everyone
      *kan komen.*                                          (HOLL_1.sml)
      can  come
      '... and properly plan something amply beforehand, so everyone can come.'

(23)  *uit   de tijd  dat  de marechaussee  nog wel*   **iets**       **beter**
      from  the  time  that  the  military-police  still  PARTICLE  something  better
      *te doen had*                                         (#dutc_4.sml)
      to  do    had
      '... from the time that the military police still had something better to do.'

(24)  *Deze Corolla is ook iets      meer geëvolueerd dan  de Subaru.*
      This  Corolla  is  also  something  more  evolved    than  the  Subaru.
      *Is net* **iets**       **beter.**                                (nie_s11.txt)
      Is  just  something  better
      'This Corolla is also a little more advanced than the Subaru. It's a just a little better.'

As for the color adjectives, the contaminating construction contains the color used as a noun, rather than as an adjective. Color adjectives are morphologically indistinguishable from color nouns in Dutch; for instance, *oranje* 'orange' can be used both as an adjective and as a noun. Accordingly, a structural ambiguity may arise in the presence of a quantifier that can be used either independently as a head or as a dependent of the noun, such as *veel* 'many/much' in (25) and (26), respectively.[21] The example in (27), *veel oranje* 'a lot of orange' is such a case of syntactical ambiguity. It can either be read as a *Dependent-Quantifier +  Head-Noun*, or as *Head-Quantifier + Partitive-Adjective*. The former syntactic structure is superficially similar to the latter, and as argued in Van de Velde (2001) and Pijpops and Van de Velde (2014), this similar construction contaminates the partitive genitive structure by rubbing off its categorically -*s*-less morphology onto the partitive color adjectives.

Interestingly, this contamination even seems to occur in unambiguous instances, where the quantifier cannot be interpreted as a dependent of the noun. As (28)–(29) show, the quantifier *iets* can only be used independently as head, yet when color adjectives combine with *iets*, they still seem to show considerably less -*s* realization than other adjectives.

(25)  *Hij  drinkt  veel.*
       he    drinks  much
       'He drinks a lot.'

(26)  *Hij  drinkt  veel    wijn.*
       he    drinks  much  wine
       'He drinks a lot of wine.'

(27)  *veel    oranje*
       much  orange
       'a lot of orange'

(28)  *Hij  drinkt  iets.*
       He    drinks  something
       'He is drinking something.'

(29)  *\*Hij  drinkt  iets       wijn.*
       he    drinks  something  wine

---

**21** We are not concerned here with the theoretical discussion whether the quantifier is best seen as a modifier or a determiner if it precedes a noun. We consider both as "dependents".

In contrast to *Type-Adjective*, the new measures *Partial String Resemblance*, *String Resemblance*, and *Semantic String Resemblance* are numeric variables that aim to directly quantify how likely an instance is to be affected by a contaminating construction. For these measures, there is no intermediary step in which the entire dataset needs to be manually inspected in order to get a sense of which constructions may be contaminating the target construction under scrutiny. While we present these measures as they are applied on the Dutch partitive genitive, they could, mutatis mutandis, be applied to any target construction. *Partial String Resemblance* requires only a partial overlap in the superficial strings of the target construction and the contaminating construction. *String Resemblance* requires a full overlap between the strings, and *Semantic String Resemblance* additionally requires a form of superficial semantic resemblance.

*Partial String Resemblance* is meant to measure how often a given adjective appears without an *-s* ending in any possibly contaminating construction. Whether or not these particular instances all actually resemble partitive genitives is not taken into account. The mere number of *-s*-less occurrences outside the partitive genitive is taken to be predictive of the rate of *-s* omission in the partitive genitive. This is, of course, a coarse measure. The reasoning is that the number of *-s*-less occurrences of an adjective will be stochastically taken into account by language users when they have to use a partitive construction comprising that adjective. Adjectives that are vastly more frequent in a construction other than the partitive genitive, even when no structural ambiguity arises, will be more likely to take the form they have in that other construction. We excluded prenominal attributive positions, though, as in Dutch, this position has a competing inflection (schwa-ending).[22]

To measure *Partial String Resemblance*, we turned to the *Corpus of Spoken Dutch* (again, as the *ConDiv* corpus is not part-of-speech annotated or syntactically parsed). For each adjective, we counted the number of times it appeared in this corpus without the *-s* ending in non-attributive position, as in (30) (compare to (2)). However, as highly frequent adjectives will necessarily appear more often in such positions, our measure *Partial String Resemblance* is not based on raw counts, but is expressed in terms of a ratio; that is, the number of times the adjective appears without the *-s* ending in a non-attributive position divided by the sum total of this number and the number of times it appears in the partitive genitive construction in the *Corpus of Spoken Dutch*. This yields a measure between 0 and 1.

---

**22** There are syntactic contexts in which attributive adjectives remain uninflected, but these are considerably outnumbered by inflection, and there is a diachronic pressure to generalize the inflection (see Van de Velde and Weerman 2014).

(30) *Uh nee nu   doe 'k 't weer  helemaal    verkeerd.* (CGN, fn000308.pos)
    uh  no  now  do  I  it  again  completely  wrongly
    'Ugh, now I'm doing it in completely the wrong way again.'

As for *String Resemblance*, what we measure is how often a sequence of quantifier and adjective, i. e. a phrase, appears without an *-s* in a sentence which resembles a partitive genitive, even if it can clearly be distinguished as a noninstance of a partitive genitive on semantic grounds. This measure is purely based on superficial formal resemblance with the entire partitive genitive phrase, i. e., no semantic resemblance is needed. It is operationalized by counting for each phrase, the number of times it appeared in the *ConDiv* corpus without the *-s* ending, while not forming a partitive genitive, as in (31) (compare to (2)). For the same reason as above, this *String Resemblance* was expressed as a ratio, by dividing the number of times the phrase appeared without the *-s* ending by the sum total of this number and the number of times it appeared as a partitive genitive in the *ConDiv* corpus.

(31) *Jongeren     beseffen wel dat  er      heel    **wat**        **verkeerd***
    youngsters  realize   PART that there  whole  something  wrongly
    *loopt,*                                                    (gva1.txt)
    runs
    'Young people do realize that a whole lot is going wrong,...'

Finally, what we measure under *Semantic String Resemblance* is how often a phrase occurs in a string that is syntactically ambiguous between a partitive genitive reading and another reading. This is a highly sensitive measure that takes into account semantics as well. This measure is operationalized by counting for each phrase, the number of times it was attested in an occurrence that will eventually be thrown out of the dataset in Subsection 5.2; that is, the number of times the phrase appeared in an occurrence that is ambiguous between a partitive genitive and another construction, as in (20) (compare to (1)). Again, *Semantic String Resemblance* was calculated as a ratio by dividing the number of times the phrase appeared in an ambiguous occurrence by the sum total of this number and the number of times it appeared as an unambiguous partitive genitive.[23]

---

[23] Note that this is not a circular measure, firstly because whether or not an occurrence was judged to be ambiguous is independent of it appearing with or without *-s* ending (see Subsection 5.2), and secondly because the constructional contamination measures (i)–(iii) will be evaluated against the restricted dataset containing the unambiguous partitive genitives only.

## 5.2 Evaluation

While *Type-Adjective* was already tested on the entire dataset in Pijpops and Van de Velde (2014), it has not yet been evaluated on a dataset containing only strictly unambiguous partitive genitives. In what follows, the ambiguous partitive genitives are removed from the dataset, and *Type-Adjective*, as well as the new measures of constructional contamination, are tested against this restricted dataset.

Distinguishing between a partitive genitive and one of the contaminating constructions can be difficult for ordinary language users, leading them to mix up the two constructions, but linguists may have a hard time as well. Even though the current dataset was already manually checked to exclude all non-partitive genitive hits, some ambiguous occurrences remain, as in (20), (32), and (33). The structural ambiguities in these examples have already been introduced above.

(32)  *Programma's in de amusementssector maar er    is nog* **niets**
      programs    in the amusement_sector but   there is still nothing
      ***concreet*.**
      concrete
      'Programs in the entertainment industry, but so far, there is nothing concrete.' **[partitive genitive]**
      'Programs in the entertainment industry, but so far, nothing is concrete.' **[predicative]**

(33)  ***veel   wit,*** *geïnspireerd op sportthema's*            (DS961102.txt)
      much   whit*e* inspired     on sport_themes
      'a lot of white things, inspired on sporting themes' **[partitive genitive]**
      'a lot of the color white, inspired on sporting themes' **[color noun]**

This means that some of the occurrences in our present dataset might not actually be partitive genitives, but rather adverbial (20), predicative (32), or color noun (33) constructions. The predicative adjectives, adverbs, and color nouns appear "ex officio" without *-s* ending in these constructions, and we suspect that such ambiguous instances are most frequent with the color and assessment adjectives. As such, it is possible that the marked preference of the color and assessment adjectives for the $[+\varnothing]$ variant in Pijpops and Van de Velde (2014) is only due to the inclusion of these ambiguous occurrences in our dataset, and that the actual unambiguous partitive genitives remain unaffected by constructional contamination.

To check for this, we again manually went through the entire dataset – including the [+s] occurrences – to code each occurrence as ambiguous or unambiguous. It was important to also check the [+s] occurrences, even though the appearance of the *-s* ending in practice disambiguates these occurrences as a partitive genitives. If we were to check the [+∅] occurrences only, the occurrence in (34) would be marked *unambiguous*, while it would have been marked *ambiguous* if it had appeared in the [+∅] variant (compare (20)). This would bias our dataset against the [+∅] variant.

(34) *Ik  moet  wel  **iets**      **verkeerd-s** gedaan  hebben.*     (v_comp9.sml)
    I    must  PART something wrong-GEN done    have
    'I must have done something wrong.'

To overcome this problem, we blinded the dataset for *-s* occurrence. That is, before the manual annotation, we removed all *-s* endings from the dataset, and then assessed each occurrence as if it were a [+∅] occurrence. In the ambiguity judgments, context and semantics were taken into account as possible (dis-)ambiguating factors. Each author evaluated 2008 hits, of which 1,000 were assessed by both authors. For these overlapping occurrences, the kappa statistic was calculated, which quantifies the strength of rater agreement for categorical data. This yielded a kappa of 0.628, indicating substantial agreement between the assessors (Landis and Koch 1977: 165). Next, the occurrences which had received conflicting evaluations were discussed, and both assessors agreed on a protocol to decide on the status of each occurrence. The evaluation of non-overlapping hits were also adapted to this protocol. Finally, all ambiguous occurrences were removed from the dataset. In this way, the restricted dataset still contained 2,700 occurrences, of which 2,276 of the [+s] variant and 424 of the [+∅] variant.

Because it was already established in Pijpops & Van de Velde (2014) that the occurrence of the *-s* ending is multifactorially determined (see Table 1), it would be ill-guided to assess the predictive performance of the four measures of constructional contamination in a bivariate test. This would run the risk that the influence of these measures would be masked or exaggerated by the influence of any of the other variables in Table 1. For instance, it might be the case that those partitive genitives with a high value for *String Resemblance* also happen to be more frequently used by Belgians or in informal language, causing them to appear more often in the [+∅] variant. In order to safeguard against this, we will evaluate how the measures perform in a regression model together with all other variables which we already know to influence *-s* omission. Concretely, we will adopt the regression model in Pijpops and Van de Velde

(2014: 18), and refit it on the restricted dataset. This then yields four regression models: one for each measure.[24]

If the effect of the variable *Type-Adjective* was only to be attributed to the ambiguous occurrences, removal of these occurrences should nullify its influence. As can be seen in Figure 3, however, a random forest analysis reveals that *Type-Adjective* is in fact still by far the most important variable in its regression model (Strobl et al. 2008). The influence of *Type-Adjective* is thus not due to the ambiguous occurrences; unambiguous partitive genitives are affected just as well.



**Figure 3:** Even fitted on strictly unambiguous data, *Type-Adjective* remains by far the most important variable in its regression model (Strobl et al. 2007).

Figure 4 shows the effect plots of each of the four measures, providing a reader-friendly visualization of their influence in their respective regression models. On the y-axis, the effect plots show the estimated probabilities for the $[+\varnothing]$ variant, while keeping the other variables, i. e., *Variety, Register, Quantifier,* and *Frequency,* constant. This means that they visualize the influence of *Type-Adjective, Partial String Resemblance, String Resemblance*, and *Semantic String*

---

**24** Apart from the measure at issue, these regression models contained the fixed effects *Variety* (Flanders, the Netherlands), *Register* (chat, e-mail, mass newspaper, quality newspaper), *Quantifier* (*iets, niets, veel, wat, weinig, zoveel*), and *Frequency* (log-transformed frequency of the partitive genitive phrase) as well as the random effect *Phrase* with the individual partitive genitive phrases as separate levels. The present regression models thus no longer contain an interaction between the variables *Variety* and *Quantifier* (cf. Pijpops and Van de Velde 2014: 18). The reason for this is that the dataset no longer contains any Netherlandic occurrences of *zoveel* 'so many' in the $[+\varnothing]$ variant. As such, including the interaction would have led to problems when calculating the estimates. For the analysis, we made use of the R software (R Core Team 2014), and of the following packages in R: MASS (Venables and Ripley 2002), rms (Harrell 2013), lme4 (Bates et al. 2013), effects (Fox 2003), dplyr (Wickham and Francois 2015), party (Hothorn et al. 2006; Strobl et al. 2007, 2008), and extrafont (Chang 2014).

*Resemblance* when taking the influence of the other variables into account. The error bars indicate 95 % confidence intervals.

Figure 4 shows that *Type-Adjective* still shows a strong preference for the [+∅] variant for the assessment adjectives, even if we remove all potentially ambiguous instances from the dataset. For the color adjectives, however, the confidence interval has increased to the point where we can no longer make any reasonable claims about them. This is simply because nearly all occurrences of the color adjectives were ambiguous between partitive genitives and color noun constructions, as in (33). Removing the ambiguous occurrences only left 10 unambiguous instances of partitive genitives with color adjectives. Even so, 6 of these exhibit the [+∅] variant. As for the other measures, Figure 4 shows that they all produce the expected effect: an increased probability of *-s* omission as their value rises.[25]
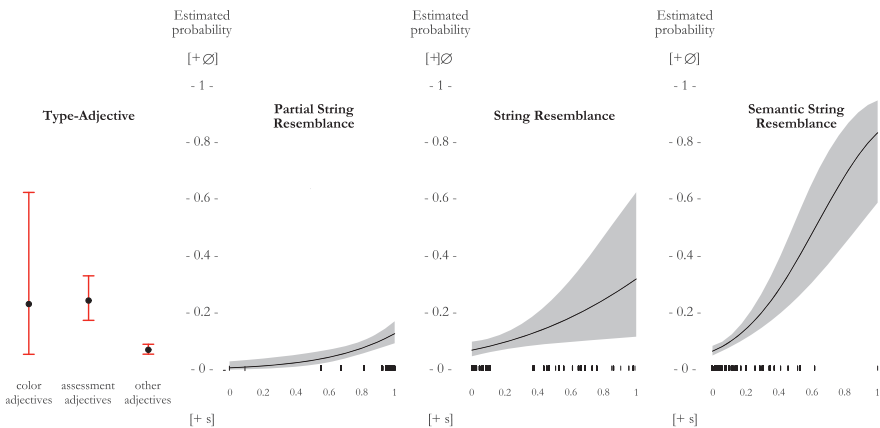


**Figure 4:** Estimated probabilities of the [+∅] variant for each of the four measures of constructional contamination, in their respective regression models. The three new resemblance measures all exhibit the expected effect: a rise in [+∅] probability as they value rises.

Let us then finally turn to a comparison of the performance of the regression models. This is done in Table 3, based on the model's AIC-value. The AIC (Akaike Information Criterion) is a measure of the unexplained variance remaining in the model, penalized for the model's complexity, and allows comparison across regression models. The lower the AIC, the better the model. Table 3 shows the AIC of each

---

**25** The "bar codes" on the x-axis, commonly known as the "rug", show the distribution of the adjectives/phrases, with each line standing for one adjective/phrase.

**Table 3:** Comparison of the regression models based on only the strictly unambiguous occurrences. *Type-Adjective* and *Semantic String Resemblance* are tied for the first place.

| Regression models | AIC | AIC-decrease due to added fixed effect | *p*-value of added fixed effect |
|---|---|---|---|
| Basic model | 1815 | | |
|   Fixed effects: *Variety, Register, Quantifier, Frequency* | | | |
|   Random effect: *Phrase* | | | |
| Basic model + Fixed effect: *Type-Adjective* | 1793 | 22 | <0.0001 |
| Basic model + Fixed effect: *Partial String Resemblance* | 1803 | 15 | 0.0002 |
| Basic model + Fixed effect: *String Resemblance* | 1811 | 4 | 0.0159 |
| Basic model + Fixed effect: *Semantic String Resemblance* | 1793 | 22 | <0.0001 |

model, as well as how much the AIC of the basic model, i. e., the model which does not contain any operationalization of constructional contamination, is decreased by adding each of the measures of constructional contamination as fixed effects. Finally, it shows whether the added variable makes a significant contribution at reducing the variance in the model, which is the case for all variables.

Of the numeric variables, *Semantic String Resemblance* is found to be the best predictor of *-s* omission, and the only one to rival the categorical variable *Type-Adjective*. What is perhaps most striking however, is the bad performance of *String Resemblance*. It barely reduces the amount of unexplained variance in the model. As such, it may be a fruitful undertaking to look under the hood of the numeric variables, to try to find out which information they encode and why some perform so much better than others. This is what is done in the next Subsection.

## 5.3 A closer look at the measures of constructional contamination

In order to get a grasp on the information encoded in the three new measures of constructional contamination, Figure 5 shows their box plots.[26] Here, it can be seen for which partitive genitive adjectives or phrases the measures exhibit extreme values.

---

26 Not all adjectives/phrases are labeled in the box plots in order to retain surveyability of the graphs.
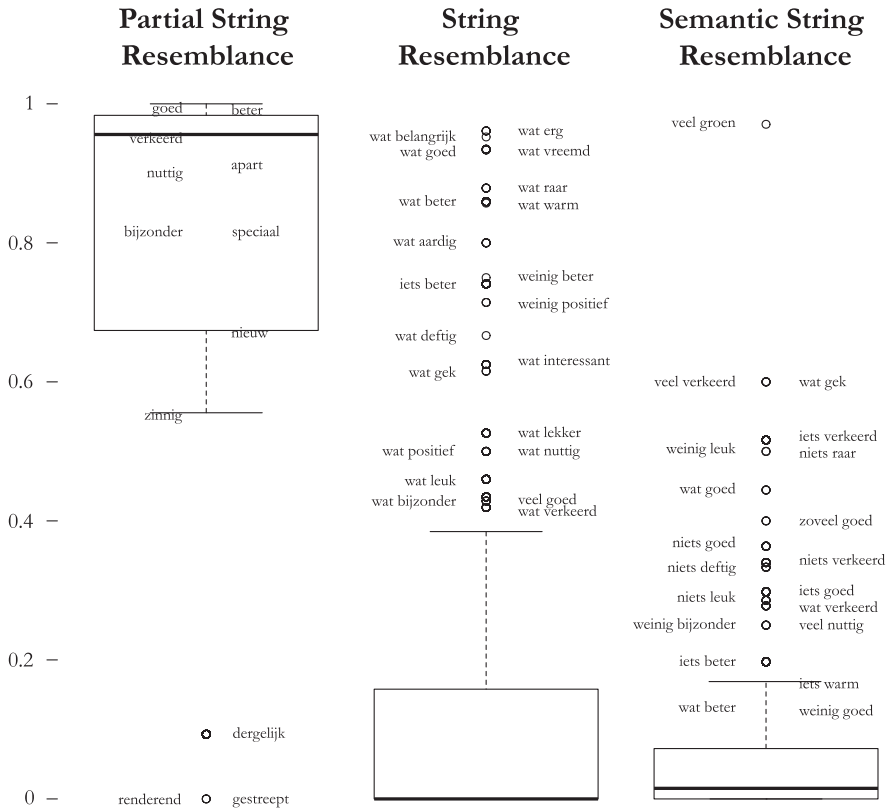
**Figure 5:** Box plots of the three new measures of constructional contamination, indicating for which adjectives or phrases the measures yield extreme values.

*Partial String Resemblance* turns out to be too coarse a measure. Most adjectives are far more frequent in non-partitive genitive contexts than in partitive genitive contexts, resulting in most values ranging between 0.95 and 1, with only tiny differences between them.[27] What this measurement mostly seems to capture is which adjectives can be found at the opposite end of the spectrum of constructional contamination, that is, which adjectives are least affected by constructional contamination, simply because they hardly appear in any (possibly contaminating) construction other than the partitive genitive.

---

[27] We tried applying an angular transformation, taking the arcsine of the square root of the values, to spread out the high values. This did not lead to a lower AIC, however.

Examples of such adjectives are *dergelijk* 'similar' and *zinnig* 'sensible'.[28] Such adjectives seem to occur so often in the partitive genitive construction that their forms with -s ending may become extremely entrenched in the minds of the language users – perhaps even more so than their 'normal' forms without -s ending – to the point that the $[+\varnothing]$ variant has little hope of catching on. As can be seen in Table 2, these adjectives do show an outspoken preference for the $[+s]$ variant.

*String Resemblance* shows, as expected, high values for the phrases with assessment adjectives. In addition, high values can also be observed for *wat erg* 'something awful', *wat vreemd* 'something strange', *wat raar* 'something weird', etc., which do not exhibit a preference for -s omission (see Table 2). In tracing the origin of these high values, we find a large number of occurrences like (35) and (36).[29] According to the reasoning behind *String Resemblance*, we would indeed expect such occurrences to generate a contaminating effect. We suspect that they do not because resemblance between these occurrences and partitive genitives is completely absent on semantic grounds. That is, their meaning is so different from partitive genitives, that no association between them and partitive genitives as in (37) and (38) is made.[30]

---

**28** The outliers *renderend* 'profitable' en *gestreept* 'striped' may be disregarded as they are extremely infrequent: they have value 0 for Adjectival Resemblance because they appeared 0 times in *The Corpus of Spoken Dutch*, and each only stands for a single data point in our dataset. They only made it into our dataset in the same way as *fout* 'incorrect' (see Subsection 4.2).

**29** This is not the case for *wat belangrijk* 'something important'. For this phrase, the high value for *String Resemblance* originates from a series of occurrences of a predicative construction containing the sequence *wat belangrijk*. However, as a partitive genitive, *wat belangrijk* is too infrequent to ascertain whether these occurrences generate a contaminating effect: it has only a single occurrence in the unambiguous dataset, which happens to exhibit the $[+s]$ variant.

**30** One could claim that formal resemblance is lacking as well, because occurrences of *[wat + ADJ]*, as in (36) and (37), mostly form complete utterances of their own, in the form of interjections, while partitive genitives are almost exclusively embedded in larger syntactic structures. Yet another possible explanation is that occurrences such as (36)–(37) are typical of Netherlandic Dutch, and therefore do not generate a contaminating effect. In Flanders, the $[+s]$ variant may have been non-native all along, introduced as a change "from above", i. e. from the Dutch standard language which is largely based on the northern varieties. It might be argued in this respect that constructional contamination only occurs when language is processed and produced in a non-native way. For instance, Flemish language users only have a vague idea what the -s ending is marking in the (northern) standard language and may therefore associate partitive genitives as in (40) more strongly with adverbial constructions as in (39) than with other partitive genitives as in (38). That would mean that constructional contamination only happens in the south. This does not seem to be the case, however: *Type-Adjective* shows the same partial effect in the south as in the north (Pijpops and Van de Velde 2014: 16–22; Pijpops and Van de Velde 2016: 359–363) and so does *Semantic String Resemblance*.

(35)  *<Vlooi> doen we tenslotte ook  voor jullie;-]* [...]
      <Vlooi> do    we after_all also for   you-PL
      *<Klubbhead>* **wat aardig***!!!!*                    (#CAIW_1.sml)
      <Klubbhead> how  nice
      '<Vlooi> After all, we're doing the same thing for you. [...] <Klubbhead>
      how nice!'

(36)  **wat erg**   *voor je:)*                             (#HASS_1.sml)
      how awful for   you
      'How awful for you.'

(37)  *Zeg ik eens* **wat       aardig** *tegen Coolgirl in een*
      say I  once something nice    to     Coolgirl in an
      *hartverwarmend DCC*                                  (#caiw_5.sml)
      heartwarming  DCC
      'For once, I say something nice to Coolgirl in an heartwarming Direct Client
      Connection,....'

(38)  *jee, heb je* **wat       erg-s**    *meegemaakt?*     (#holl_6.sml)
      gee have you something aweful-GEN experienced?
      'Gee, have you gone through something awful?'

If we supplement the formal resemblance of the phrase with semantic informa-
tion, as in the variable *Semantic String Resemblance*, we find that performance
markedly improves. At this point, semantic resemblance is a categorical dis-
tinction: either the occurrence is unequivocally not a partitive genitive, and
then there is no semantic resemblance, or the occurrence is semantically
ambiguous between a reading as a partitive genitive and some other construc-
tion, and then there is semantic resemblance. Ideally, we would want to say
that a sentence such as (2), repeated below as (39), is semantically more
resembling of a partitive genitive than (35)–(36), even though in both cases
the context makes clear that they are not partitive genitives. Still, we have the
intuition that 'interpret something in a wrong way' (39) is semantically more
resembling of 'say something which is wrong' (40), than 'how nice!' (35) is
resembling of 'something which is nice' (37). However, we are not aware of a
way to objectively quantify this intuition, at least without carrying out exten-
sive surveys with native speakers, which falls beyond the scope of the present
study.

(39)  *dat* **iets** **verkeerd** *geïnterpreteerd wordt?*     (#VLAAN_1.sml)
      that something wrongly  interpreted    gets
      '...that something gets wrongly interpreted?'

(40)  *Oi  nu   de   bek   houden voor   ik* **iets** **verkeerd** *zeg*
      oi  now  the  beak  hold    before I  something wrong    say
                                                      (n_comm6.sml)
      'Oi, I'd better shut my mouth before I say something wrong.'

The final measure, *Semantic String Resemblance*, thus produces the best results. It seems both formal and semantic resemblance between constructions is needed to trigger constructional contamination. We can now further develop the diagram in Figure 1, resulting in the one in Figure 6. The corpus examples which are referred to in Figure 6 are repeated below, for the comfort of the reader.
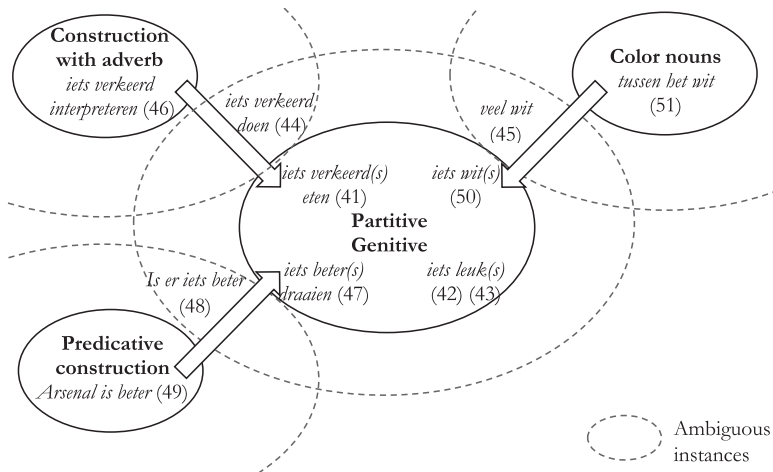


**Figure 6:** Result of applying Figure 1 to the partitive genitive. Constructional contamination affecting the partitive genitive through superficial resemblances, with corpus examples.[31]

---

**31** Note that the corpus examples in this diagram are nothing more than just examples. For instance, contamination from the adverbial construction affects many more adjectives than *verkeerd* 'wrong'. In fact, the adjective *beter* 'better', which is used here as an example of the predicative construction, also very often appeared in adverbial occurrences resembling partitive genitives.

(41) *in begin     van de week **iets**     **verkeerd** gegeten* (#LEUV_4.sml)
     in beginning of   the week something wrong     eaten
     'I ate something wrong at the start of the weak.'

(42) *Is er     nog **iets**     **leuk-s** te beleven?*          (#holl_6.sml)
     is there still something fun-GEN to experience?
     'Is there still something fun to do?'

(43) *of        er     hier   nog **iets**         **leuk** te beleven     valt*
     whether there here still something fun  to experience falls
                                                              (#holl_6.sml)
     '... whether there is still something fun to do here?'

(44) *Ze    zijn dus vaak erg  onzeker  en    bang  dat  ze*
     they  are  thus often very insecure and  afraid that they
     ***iets**         **verkeerd** zullen doen.*         (verstr4.txt)
     something wrong       will    do
     'As such, they are often very insecure and afraid to do something wrong.'
     **[partitive genitive]**
     'As such, they are often very insecure and afraid to do something wrongly.'
     **[construction with adverb]**

(45) ***veel    wit,**  *geïnspireerd op sportthema's*          (DS961102.txt)
     much  whit*e* inspired       on sport_themes
     'a lot of white things, inspired on sporting themes' **[partitive genitive]**
     'a lot of the color white, inspired on sporting themes' **[color noun]**

(46) *dat iets     **verkeerd** geïnterpreteerd wordt?*       (#VLAAN_1.sml)
     *that something wrongly   interpreted     gets*
     *'...that something gets wrongly interpreted?'*

(47) *Zal  wel  backstreet boys zijn of zo  ... **iets***
     will  PART backstreet boys be   or_so    something
     ***beter-s**  *draaien ze    daar toch   niet!*          (#NEDE_1.sml)
     better-GEN play     they there PARTICLE not
     'It will probably be the Backstreet Boys or something like that. They don't
     play anything better there.'

(48)   *Is er    dan   **iets**      **beter**  dan  een  SB Live?*      (#vlaan_8.sml)
        is there  then  something  better  than  an   SB  Live?
        'Does something better than an SB Live exist, then?' **[partitive genitive]**
        'Does something then exist which is better than an SB Live?' **[predicative construction]**

(49)   *Ar*senal  *is*  ***beter.***                  (CGN, fn007885.pos)
        Arsenal  is  better
        'Arsenal is better.'

(50)   *<Tone> Sneeuw? <Tone> Wasda?*       *<bruintje> ja,  da*
        <Tone> snow     <Tone>  what = is = that  <bruintje>  yes  that
        *is  **iets**      **wit.***                 (#DUTC_2.sml)
        is  something  white
        '<Tone> Snow? <Tone> What's that? <bruintje> Well yeah, that's something white...'

(51)   *Awel tussen   **het wit**   heb  je    een zwart*       (#DUTC_2.sml)
        well   between  the  white  have  you  a     black
        'Well, between the white, you have a black...'

# 6 Conclusions

Constructional contamination has been defined in Sections 1 and 2. We can now summarize the answer to the question posed in the title of this paper: how does constructional contamination work and how do we measure it?

Constructional contamination works through superficial formal and semantic resemblance. The required formal resemblance is only superficial in that it involves strings that are similar at the surface level, not at any deeper syntactic level. Meanwhile, the required semantic resemblance is only superficial in that the meaning of both constructions can in fact be quite different, as long as they are similar enough to allow for the existence of some ambiguous instances.

For contamination to take place, a bridge needs to be established between the contaminating and the target construction. This bridge is formed by ambiguous occurrences in which both readings are possible. Once this bridge is in place, constructional contamination may affect even the very core of the

contaminated construction, i.e., completely unambiguous occurrences of the target construction. In order to speak of constructional contamination, it is of paramount importance that these occurrences are affected as well. The observation that an ambiguous occurrence is affected by one of the constructions at issue is trivial; this is probably because it simply *is* an instance of this construction (see Section 5).

For the partitive genitive in Dutch, constructional contamination can be observed in that the frequent occurrence of the string *iets verkeerd* 'something wrong', and comparable examples in sentences that superficially resemble partitive genitives, affects the realization of the genuine, unambiguous partitive genitive phrase *iets verkeerd(s)* 'something wrong'. Strings that do not occur in the contaminating construction, however, remain unaffected.

We believe that the mechanism of constructional contamination, as measured in this study, bears testimony to the hypothesis that language users do not always analyze sentences such as (1)–(2) and (39)–(41) to their underlying structures, but instead only chunk them into large, unanalyzed exemplars such as *iets verkeerd,* which are stored and accessed as wholes, especially when these chunks are frequent (see Bybee 2010 and Diessel 2007 for the pervasive effects of frequency). This would allow the frequent recurrence of sentences like (2) to affect the realization of partitive genitives. There is corroborating evidence for the template nature of linguistic constructions (Dąbrowska 2014), with speakers relying on a "quick-and-dirty pseudo-parse", which underspecifies the syntactic structure of utterances (Ferreira and Patson 2007, cited in Dąbrowska 2012). One may take this a step further, as Bauer does (1983: 296, quoted in Hüning 1999: 30): "It might (...) be worth speculating whether language users work by analogy whereas linguists interpret such behavior in terms of rules, so that a linguist's description is inevitably a fiction."

If we shift our attention from the individual language user to the community level (Verhagen 2013; Dąbrowska 2015), our findings mesh with the idea that language as a whole is emergent, and that the constructions it is composed of have a temporary, transient or ephemeral status, much like moving sand dunes (Hopper 1987, 1998). There is no discrete boundary between constructions, and features may travel horizontally from one construction to the next, on the basis of superficial formal and semantic resemblance (De Smet and Van de Velde 2014; Norde 2014), forming the basis of multiple source constructions in diachrony (De Smet et al. 2013).

As to measuring constructional contamination, we have proposed a measure in Section 5 that takes into account both formal and semantic resemblence. This measure, called *Semantic String Resemblance*, outperformed other measures in the present case study of the Dutch partitive genitive. While, as noted in

Subsection 5.3, there still seems to be some room to further refine its operationalization, this measure can already be applied to other case studies.

As a final concluding remark, we would like to point out that this study can serve as an example of the usefulness of manually sifting through authentic language data, as opposed to both introspective judgments of self-created sentences and large-scale automatic handling of corpus data. The only way we caught track of constructional contamination was through the manual checking of our dataset. At the same time, we should not shy away from quantitative methods in linguistics. For some old-school philologists, there is no distinction between automatic annotation of corpus data and applying statistics. This is a serious misconception. In the paper at hand, the only way to show the superiority of the hand-coded semantically sensitive variable over a coarser semi-automatically coded variable was to apply statistical methods.

We would therefore like to end this paper with a call to view neither advanced statistics, nor manual annotation as necessary evils when doing corpus research. Rather, the former is a welcome tool to bolster the empirical enterprise in linguistics, and the latter is a valuable way of keeping direct contact between researcher and data.

# References

Abbot-Smith, Kirsten & Heike Behrens. 2006. How known constructions influence the acquisition of other constructions: The German passive and future constructions. *Cognitive Science: A Multidisciplinary Journal of Artificial Intelligence, Linguistics, Neuroscience, Philosophy, Psychology* 30(6). 995–1026.

Anthony, Laurence. 2011. *AntConc (Computer Software, version 3.3.3)*. Tokyo: Waseda University. http://www.antlab.sci.waseda.ac.jp/.

Arnon, Inbal & Neal Snider. 2010. More than words: Frequency effects for multi-word phrases. *Journal of Memory and Language* 62(1). 67–82.

Baayen, Rolf Harald. 2008. Analyzing linguistic data: A practical introduction to statistics using R. Cambridge: Cambridge University Press.

Balk-Smit Duyzentkunst, Frida. 2000. *Grammatica van het Nederlands* [Grammar of Dutch]. Den Haag: Sdu.

Bartol, Thomas, Cailey Bromer, Justin Kinney, Michael Chirillo, Jennifer Bourne, Kristen Harris & Terrence Sejnowski. 2015. Nanoconnectomic upper bound on the variability of synaptic plasticity. *eLife* 4. e10778.

Bates, Douglas, Martin Maechler, Ben Bolker & Steven Walker. 2013. *lme4: Linear mixed-effects models using Eigen and S4. R package version 1.0-4*. http://cran.r-project.org/package=lme4.

Bauer, Laurie. 1983. *English word-formation*. Cambridge: Cambridge University Press.

Bergs, Alexander & Lena Heine. 2010. Mood and modality in English. In Rolf Thieroff & Björn Rothstein (eds.), *Mood systems in the languages of Europe*, 103–117. Amsterdam: John Benjamins.

Booij, Geert. 2010. *Construction morphology*. Oxford: Oxford University Press.

Broekhuis, Hans. 2013. *Syntax of Dutch: Adjectives and adjective phrases*. Amsterdam: Amsterdam University Press.

Bybee, Joan. 2010. *Language, usage and cognition*. Cambridge: Cambridge University Press.

Chang, Winston. 2014. extrafont: Tools for using fonts. http://cran.r-project.org/package=extrafont.

Croft, William. 2001. *Radical construction grammar: Syntactic theory in typological perspective*. Oxford: Oxford University Press.

Dąbrowska, Ewa. 2012. Different speakers, different grammars: Individual differences in native language attainment. *Linguistic Approaches to Bilingualism* 2(3). 219–253.

Dąbrowska, Ewa. 2014. Recycling utterances: A speaker's guide to sentence processing. *Cognitive Linguistics* 25(4). 617–653.

Dąbrowska, Ewa. 2015. Language in the mind and in the community. In Jocelyne Daems, Eline Zenner, Kris Heylen & Dirk Speelman (eds.), *Change of paradigms – new paradoxes: Recontextualizing language andlLinguistics*. Berlin: De Gruyter Mouton.

De Smet, Hendrik. 2010. Grammatical interference: Subject marker *for* and the phrasal verb particles *out* and *forth*. In Elizabeth Trousdale & Graeme Traugott (eds.), *Gradience, gradualness and grammaticalization*, 75–104. Amsterdam: John Benjamins.

De Smet, Hendrik, Lobke Ghesquière & Freek Van de Velde (eds.). 2013. On multiple source constructions in language change. [Special issue] *Studies in Language* 37(3).

De Smet, Hendrik & Freek Van de Velde. 2013. Serving two masters: Form–function friction in syntactic amalgams. *Studies in Language* 37(3). 534–565.

De Smet, Hendrik & Freek Van de Velde. 2014. Travelling features: Multiple sources, multiple destinations. Paper presented at The 8th International Conference on Construction Grammar (ICCG8), University of Osnabrück, 2–6 September.

Diessel, Holger. 2007. Frequency effects in language acquisition, language use, and diachronic change. *New Ideas in Psychology* 25. 108–127.

Diessel, Holger. 2015. Usage-based construction grammar. In Ewa Dąbrowska & Dagmar Divjak (eds.), *Handbook of Cognitive Linguistics*, 296–321. Berlin: De Gruyter Mouton.

Ferreira, Fernanda & Nikole Patson. 2007. The "good enough" approach to language comprehension. *Language and Linguistics Compass* 1. 71–83.

Fonteyn, Lauren & Nikki van de Pol. 2016. Divide and conquer: The formation and functional dynamics of the Modern English *ing*-clause network. *English Language and Linguistics* 20(2). 185–219.

Fox, John. 2003. Effect displays in R for generalised linear models. *Journal of Statistical Software* 8. 1–27.

Goldberg, Adele Eva. 1995. *Constructions: A construction grammar approach to argument structure.* Chicago: University of Chicago press.

Gries, Stefan Th. 2013a. *Statistics for linguistics with R: A practical introduction*, 2nd edn. Berlin & New York: De Gruyter.

Gries, Stefan Th. 2013b. 50-something years of work on collocations: What is or should be next. *International Journal of Corpus Linguistics* 18(1). 137–165.

Gries, Stefan Th. 2014. Coll.analysis 3.5. A script for R to compute perform collostructional analyses.

Gries, Stefan Th. & Anatol Stefanowitsch. 2004. Extending collostructional analysis: A corpus-based perspective on "alternations." *International Journal of Corpus Linguistics* 9(1). 97–130.

Grondelaers, Stefan, Katrien Deygers, Hilde Van Aken, Vicky Van den Heede & Dirk Speelman. 2000. Het CONDIV-corpus geschreven Nederlands [The CONDIV-corpus of written Dutch]. *Nederlandse Taalkunde* 5(4). 356–363.

Haeseryn, Walter, Kirsten Romijn, Guido Geerts, Jaap de Rooij & Maarten van den Toorn. 1997. *Algemene Nederlandse Spraakkunst* [General Dutch Grammar]. Groningen: Nijhoff.

Halliday, Michael & Christian Matthiessen. 2004. *An introduction to functional grammar*, 3rd edn. London: London Arnold.

Harrell, Frank. 2013. rms: Regression modeling strategies. R package version 4.0-0. http://cran.r-project.org/package = rms.

Heine, Bernd. 2002. On the role of context in grammaticalization. In Ilse Wisher & Gabriele Diewald (eds.), *New reflections on grammaticalization*, 83–101. Amsterdam: John Benjamins.

Hoeksema, Jack. 2014. De opkomst van "aan" als verbindend element in maatnomenconstructies [The rise of "aan" as a connecting element in measure noun constructions]. In Freek Van de Velde, Hans Smessaert, Frank Van Eynde & Sara Verbrugge (eds.), *Patroon en argument*: *Een dubbelfeestbundel bij het emeritaat van William Van Belle en Joop van der Horst* [Pattern and argument: A double festschrift on the occasion of William Van Belle's and Joop van der Horst's retirement], 421–432. Leuven: Leuven University Press.

Hopper, Paul. 1987. Emergent grammar. *Berkeley Linguistic Society* 13. 139–157.

Hopper, Paul. 1998. Emergent grammar. *The new psychology of language: Cognitive and functional approaches to language structure*, 155–175. Mahwah, NJ: Lawrence Erlbaum.

Horst, Joop van der. 2008. *Geschiedenis van de Nederlandse syntaxis* [History of Dutch syntax]. Leuven: Universitaire Pers Leuven.

Hothorn, Torsten, Peter Bühlmann, Sandrine Dudoit, Annette Molinaro & Mark Van Der Laan. 2006. Survival ensembles. *Biostatistics* 7(3). 355–373.

Huddleston, Rodney. 2002. The verb. In Rodney Huddleston & Geoffrey Pullum (eds.), *The Cambridge grammar of the English language*, 71–212. Cambridge: Cambridge University Press.

Hüning, Matthias. 1999. *Woordensmederij: De geschiedenis van het suffix -erij* [Word forging: The history of the suffix *-erij*]. The Hague: The Hague Holland Academic Graphics.

Hunston, Susan & Geoff Thompson (eds.). 2001. *Evaluation in text: Authorial stance and the construction of discourse*. Oxford: Oxford University Press.

Jurafsky, Daniel. 1992. *An on-line computational model of human sentence interpretation: A theory of the representation and use of linguistic knowledge*. Berkeley, CA: University of California dissertation.

Lakoff, George. 1988 [1974]. Syntactic amalgams. In Eric Schiller, Barbara Need, Douglas Varley & William Eilfort (eds.), *The best of CLS*: *A selection of out-of-print papers from 1968 to 1975*, 25–45. Chicago: Chicago Linguistic Society.

Lakoff, George. 1987. *Women, fire, and dangerous things. What categories reveal about the mind*. Chicago: University of Chicago Press.

Landis, John Richard & Gary Grove Koch. 1977. The measurement of observer agreement for categorical data. *Biometrics* 33(1). 159–174.

Langacker, Ronald W. 2008. *Cognitive grammar: A basic introduction*. Oxford: Oxford University Press.

Markman, Ellen & Gwyn Wachtel. 1988. Children's use of mutual exclusivity to constrain the meanings of words. *Cognitive Psychology* 20(2). 121–157.

Markman, Ellen, Judith Wason & Mikkel Hansen. 2003. Use of the mutual exclusivity assumption by young word learners. *Cognitive Psychology* 47(3). 241–275.

Martin, James & Peter White. 2007. *The language of evaluation: Appraisal in English*. Houndmills: Palgrave Macmillan.

Norde, Muriel. 2014. On parents and peers in constructional networks. Paper presented as CoglingDays 6. University of Ghent, December 12.

Oostdijk, Nelleke, Wim Goedertier, Frank Van Eynde, Louis Boves, Jean-Pierre Martens, Michael Moortgat & Harald Baayen. 2002. Experiences from the spoken Dutch corpus project. *Proceedings of the third International Conference on Language Resources and Evaluation*, 340–347. Las Palmas. http://www.lrec-conf.org/proceedings/lrec2002/.

Pijpops, Dirk & Freek Van de Velde. 2014. A multivariate analysis of the partitive genitive in Dutch: Bringing quantitative data into a theoretical discussion. *Corpus Linguistics and Linguistic Theory*. Published online, ahead of print.

Pijpops, Dirk & Freek Van de Velde. 2016. Ethnolect speakers and Dutch partitive adjectival inflection: A corpus analysis. *Taal en Tongval* 67(2). 343–371.

R Core Team. 2014. *R: A language and environment for statistical computing. R Foundation for Statistical Computing*. Vienna. http://www.r-project.org/.

Speelman, Dirk. 2014. Logistic regression: A confirmatory technique for comparisons in corpus linguistics. In Dylan Glynn & Justyna A. Robinson (eds.), *Corpus methods for semantics: Quantitative studies in polysemy and synonymy*, 487–533. Amsterdam: John Benjamins.

Strobl, Carolin, Anne-Laure Boulesteix, Thomas Kneib, Thomas Augustin & Achim Zeileis. 2008. Conditional variable importance for random forests. *BMC Bioinformatics* 9(307). Available at http://bmcbioinformatics.biomedcentral.com/articles/10.1186/1471-2105-9-307. DOI: 10.1186/1471-2105-9-307.

Strobl, Carolin, Anne-Laure Boulesteix, Achim Zeileis & Torsten Hothorn. 2007. Bias in random forest variable importance measures: Illustrations, sources and a solution. *BMC Bioinformatics* 8(25). Available at http://bmcbioinformatics.biomedcentral.com/articles/10.1186/1471-2105-8-25. DOI: 10.1186/1471-2105-8-25.

Szmrecsanyi, Benedikt. 2005. Language users as creatures of habit: A corpus-linguistic analysis of persistence in spoken English. *Corpus Linguistics and Linguistic Theory* 1(1). 113–150.

Townsend, David & Thomas Bever. 2001. *Sentence comprehension: The integration of habits and rules*. Cambridge, MA: MIT Press.

Van Bart, Peter, Johan Kerstens & Arie Sturm. 1998. *Grammatica van het Nederlands: Een inleiding* [Grammar of Dutch: An introduction]. Amsterdam: Amsterdam University Press.

Van de Velde, Freek. 2001. *Iets taalkundig(s): Een functioneel georiënteerde analyse van deflexie en de genitiefontwikkeling in het Nederlands* [Something linguistic: A functionally oriented analysis of deflexion and the development of the genitive in Dutch]. Leuven: Univerisity of Leuven MA thesis.

Van de Velde, Freek. 2009. *De nominale constituent: Structuur en geschiedenis* [The noun phrase. Structure and history]. Leuven: Leuven University Press.

Van de Velde, Freek. 2014. Degeneracy: The maintenance of constructional networks. In Ronny Boogaart, Timothy Colleman & Gijsbert Rutten (eds.), *Extending the scope of Construction Grammar*, 141–179. Berlin: De Gruyter Mouton.

Van de Velde, Freek, Hendrik De Smet & Lobke Ghesquière. 2013. On multiple source constructions in language change. *Studies in language* 37(3). 473–489.

Van de Velde, Freek & Joop van der Horst. 2013. Homoplasy in diachronic grammar. *Language Sciences* 36(1). 66–77.

Van de Velde, Freek & Fred Weerman. 2014. The resilient nature of adjectival inflection in Dutch. In Petra Sleeman, Freek Van de Velde & Harry Perridon (eds.), *Adjectives in Germanic and Romance*, 113–145. Amsterdam: John Benjamins.

Venables, William & Brian Ripley. 2002. *Modern applied statistics with S*, 4th edn. New York: Springer.

Verhagen Arie. 2013. Darwin en de ideale taalgebruiker [Darwin and the ideal language user]. In Theo A.J.M. Janssen & Jan Noordegraaf (eds.), *Honderd jaar taalwetenschap. Artikelen aangeboden aan Saskia Daalder bij haar afscheid van de Vrije Universiteit* [A hundred years of linguistics. Articles presented to Saskia Daalder on the occasion of her retirement from the Free University], 151–162. Amsterdam/Münster: Stichting Neerlandistiek VU/ Nodus Publikationen.

Wickham, Hadley & Romain Francois. 2015. dplyr: A grammar of data manipulation. http://cran. r-project.org/package = dplyr.

Zipf, George Kingsley. 1932. *Selected studies of the principle of relative frequency in language*. Harvard: Harvard University Press.